# Plant science data management and integration: the heterogeneity and dispersion challenge
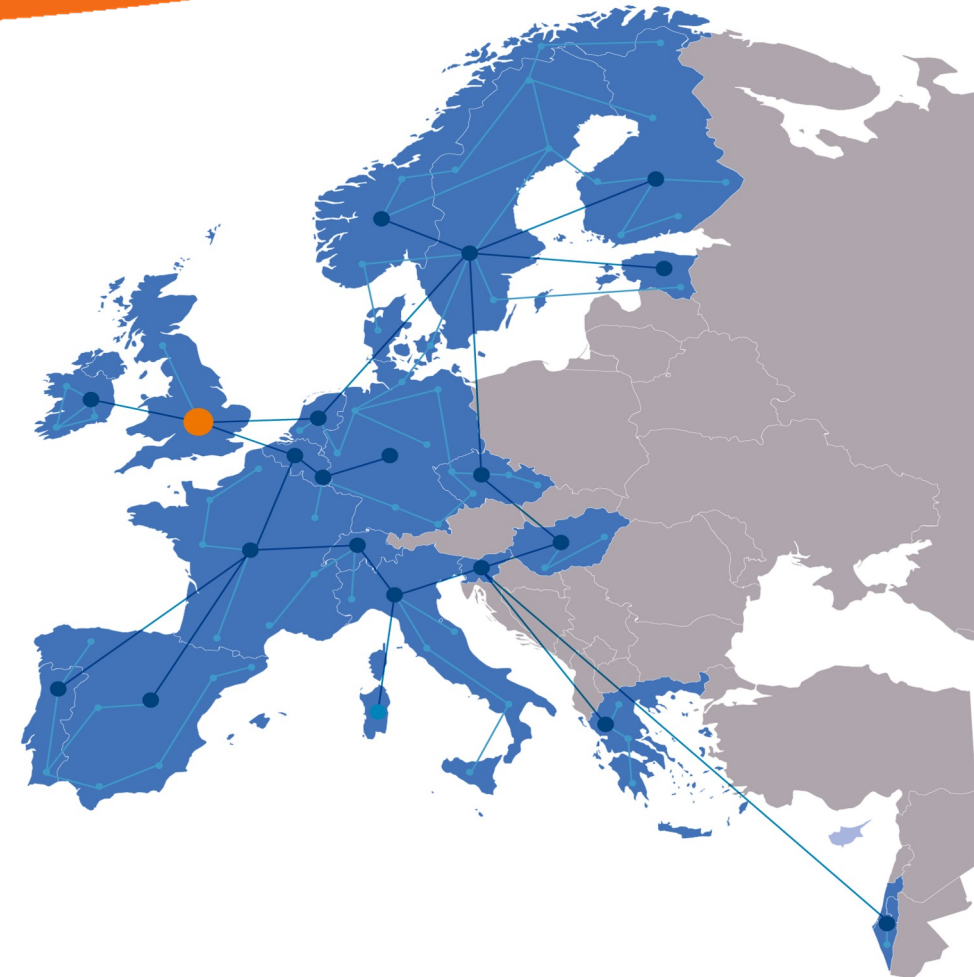## Solutions from ELIXIR and EMPHASIS European Infrastructure and beyond

*www.elixir-europe.org*

# ELIXIR

- intergovernmental organisation
- life science resources
  - databases, software tools, training materials, standards, compute resources
  - across Europe.

coordinate life science resources ➔ single

infrastructure:

- Find and share data
- Exchange expertise
- Agree on best practises in scientific research

# ELIXIR's work

**ELIXIR** coordinates activities through at least one of five areas of activities called Platforms:

- Compute
- Data
- Interoperability
- Tools
- Training

Driven by eleven **ELIXIR** Communities

**Plant Sciences Community**
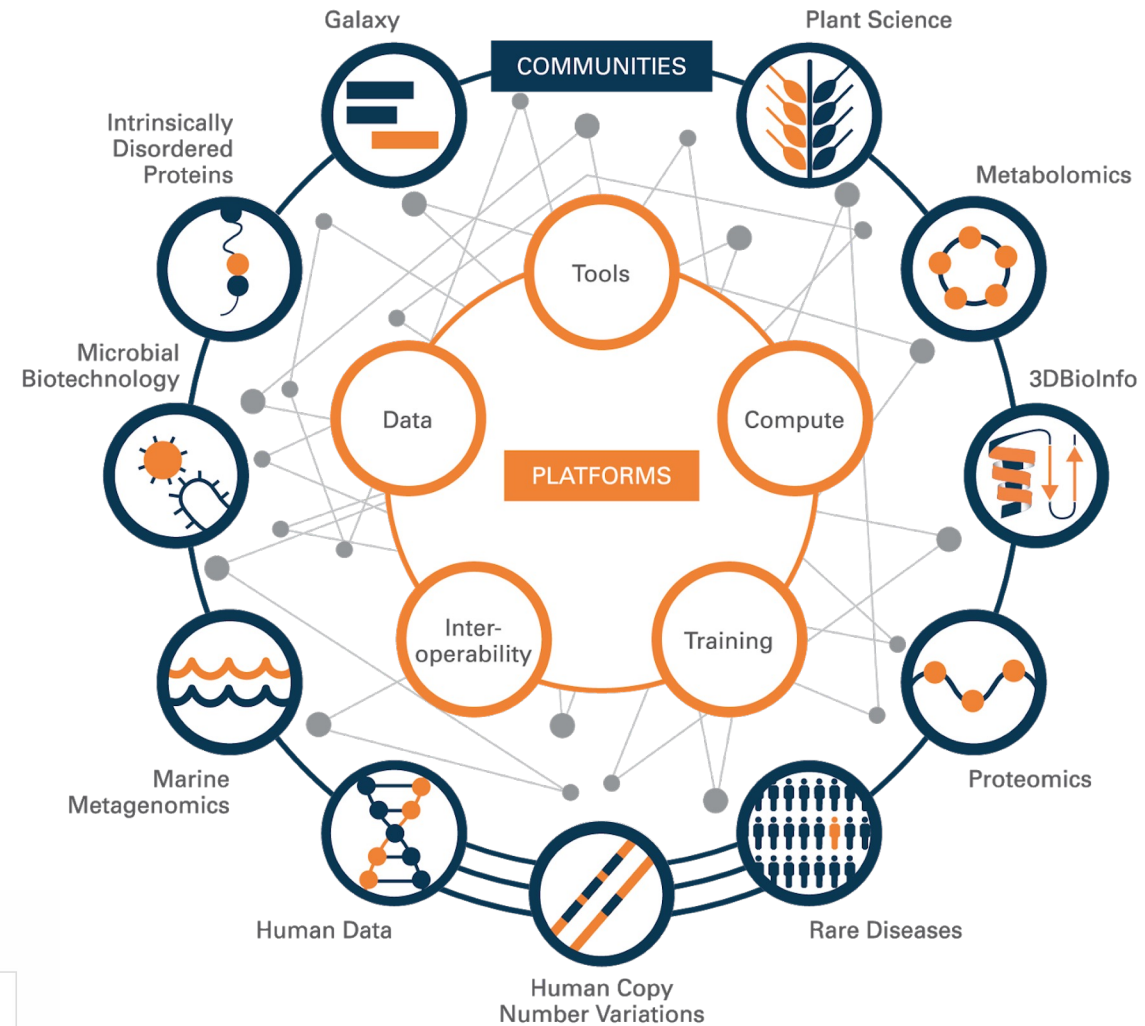


### Leadership

**Sebastian Beier**
(ELIXIR Germany)

**Kristina Gruden**
(ELIXIR Slovenia)

**Cyril Pommier**
(ELIXIR France)

**Katharina Heil**
(Communities Coordinator, ELIXIR Hub)

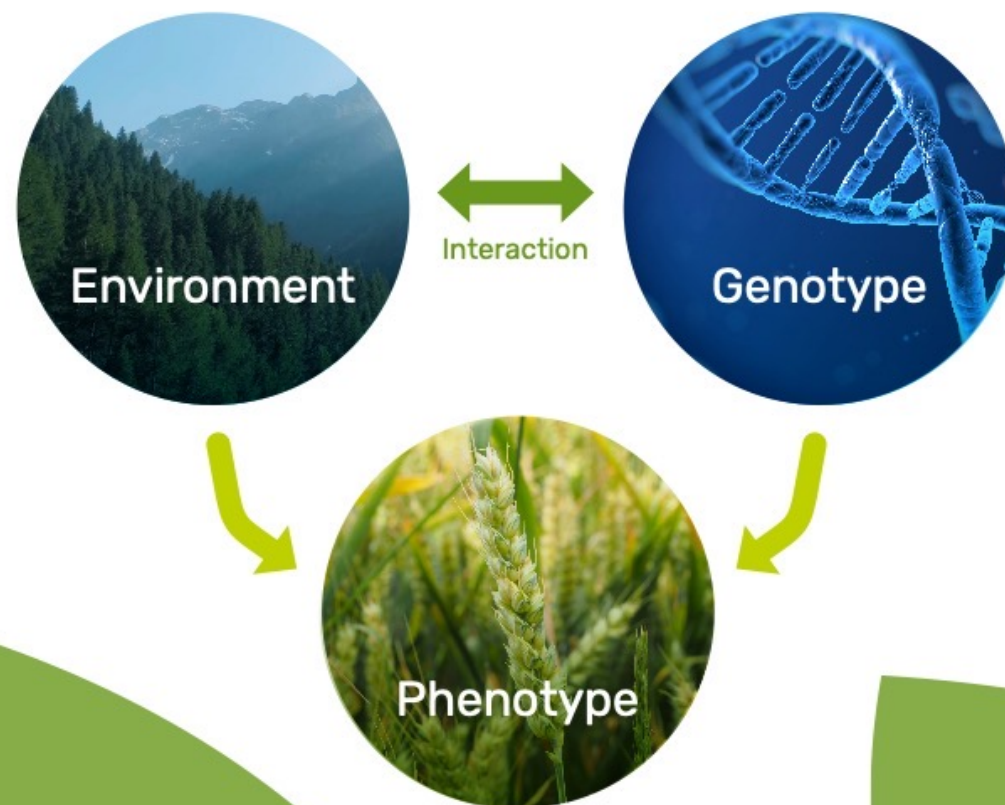https://elixir-europe.org/communities/plant-sciences

# EMPHASIS

## EUROPEAN INFRASTRUCTURE FOR PLANT PHENOTYPING

### SCIENTIFIC TOOL TO STUDY PLANT–ENVIRONMENT INTERACTION

- Study of plant structure and function

- Using non-invasive technology

- Understanding how plant structure and function depend on genetics and the environment

### How does a plant cope with its environment?

Phenotyping is used to understand how plants can cope with reduced resources, pathogens and climate change.

Environment ⟷ Interaction ⟷ Genotype

Phenotype

# Infrastructure Categories

PLANT PHENOTYPING REQUIRES INTEGRATION OF BOTH FACILITIES AND ACTIVITIES
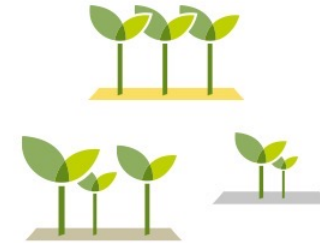
**CONTROLLED CONDITIONS**

Investigation of diverse plant traits in response to well-defined environmental conditions

**INTENSIVE FIELD**

Detailed investigation of plants and canopies under well-monitored field conditions

**LEAN FIELD**

Field sites with basic equipment and environmental monitoring that can be linked to a network of field sites

**MODELLING**

Models integrated in phenotyping pipelines and predictive models using phenotypic data

**DATA & COMPUTATIONAL SERVICES**

Integrating compatible information systems to provide access to data

elixir

# Open science through FAIR data principles

Wilkinson et al., *The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data 3 (2016)*



**Findable**

Ids

Index
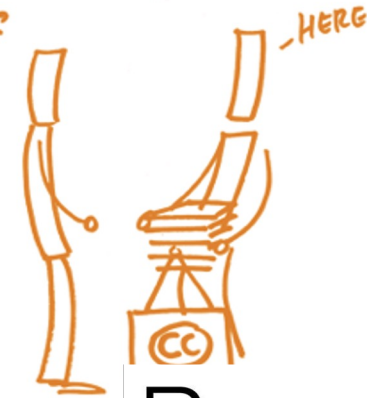
Metadata

Description

**Accessible**

Open Protocole

Perenial Metadata

**Interoperable**

Semantics

Linked Data

Vocabularies

**Reusable**

License

Well described

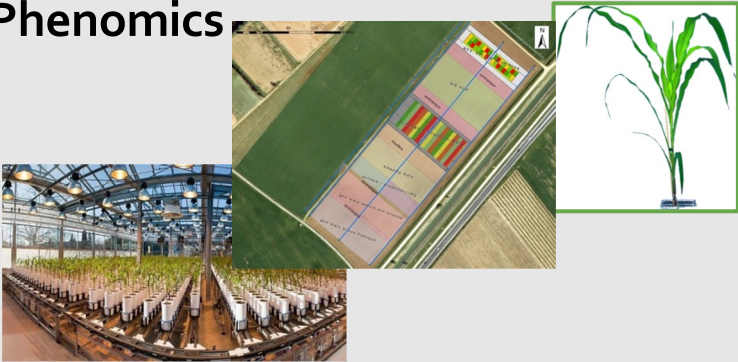Provenance (origin, process, methodology)
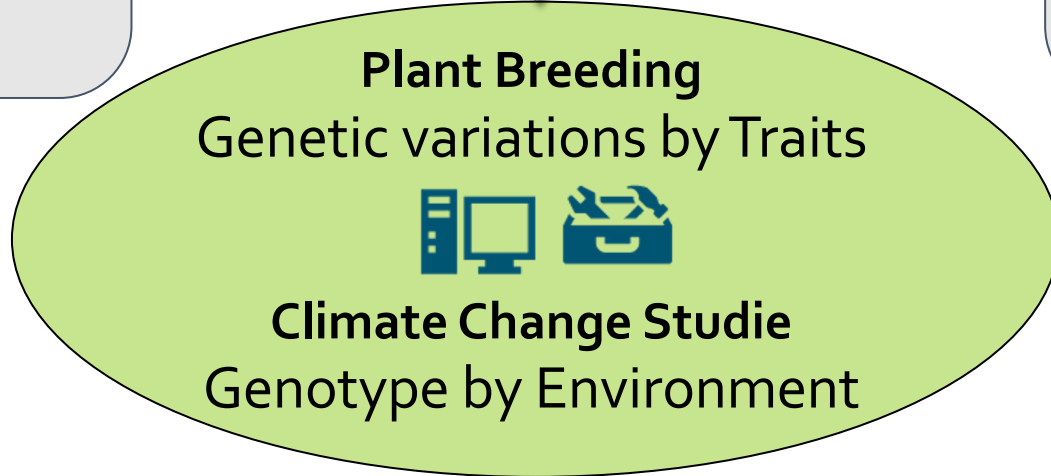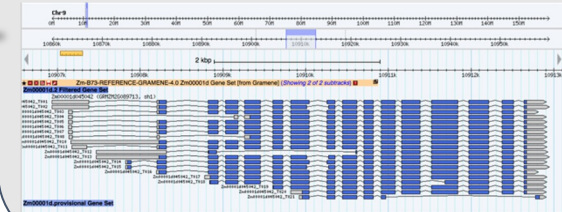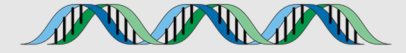
Standards

Sustainable data access over decades

INRAe

Plant science data management and integration: the heterogeneity and dispersion challenge
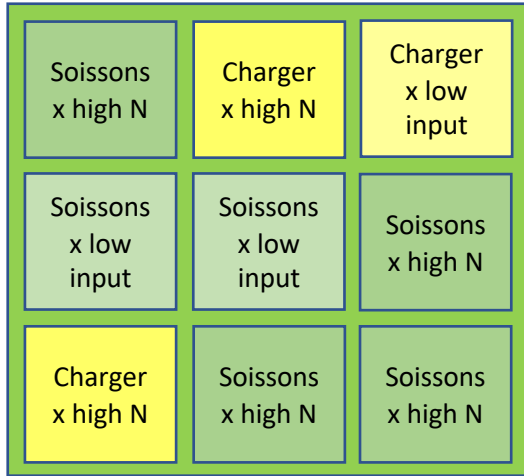11 March 2022 / ABLS4 / Cyril Pommier

p. 6

# > What is FAIR for plant data ?

- Phenotyping
  - Raw data
    - Images
    - NIRS
    - Individual plant time series
    - Expensive to generate
    - Not reproducible
  - Computed / derived data
    - Data matrices (XLSX)
- Genetic variation
  - Raw data
    - Sequence files
    - "cheap" to generate
    - Big Data
  - Derived
    - VCF
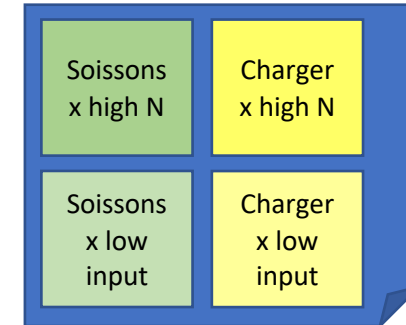    - Aligned to a given reference genome

INRA

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 9

« Raw » data, pheno/env measures, variables

« computed » data, reduced, indicators

Derivation, Reduction

| Soissons x high N | Charger x high N | Charger x low input |
| Soissons x low input | Soissons x low input | Soissons x high N |
| Charger x high N | Soissons x high N | Soissons x high N |

| Soissons x high N | Charger x high N |
| Soissons x low input | Charger x low input |

| Genotype | Treatment | N input | Date | Rep | Fusariose |
|---|---|---|---|---|---|
| Soissons | low input | 15,32253129 | 15/11/2011 | 1 | 5 |
| Soissons | low input | 15,31430556 | 16/11/2011 | 2 | 7 |

| Genotype | Treatment | Fusariose |
|---|---|---|
| Soissons | low input | 6 |

| | | | | |
|---|---|---|---|---|
| 661300270 Ardon | 2005 | | | |
| 661300444 Ardon | 2004 38.96112577281653 | 12/01/2004 228.8 | | |
| 661300444 Ardon | 2005 | | | |
| 661300312 Cavallermaggiore | 2004 52.4 | 01/01/2004 249.9 | | |
| 661300312 Cavallermaggiore | 2005 | | | |
| 661300371 Cavallermaggiore | 2004 45.74 | 01/01/2004 230.2 | | |
| 661300371 Cavallermaggiore | 2005 | | | |
| 661300487 Cavallermaggiore | 2004 72.52 | 01/01/2004 309.8 | | |
| 661300487 Cavallermaggiore | 2005 | | | |
| 661300585 Cavallermaggiore | 2004 71.739999999999995 | 01/01/2004 305.7 | | |
| 661300585 Cavallermaggiore | 2005 | | | |
| 661300468 Headley | 2004 45.27 | 01/01/2004 | | |
| 661300468 Headley | 2005 | | | |
| 661300469 Headley | 2004 70.930000000000007 | 01/01/2004 | | |
| 661300469 Headley | 2005 | | | |

INRAE

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

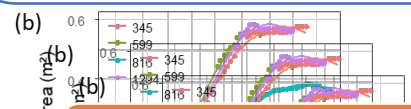p. 10

# Plant Phenotyping Life cycle

**Data acquisition**

- **VARIABLES**
- Plant/microplot level
- Traceability
- Raw measures
- Data Cleaning
- Platform IS (Emphasis IS, PHIS, …)
- Analysis Reproducibility
- Provenance

**Data computation**

- **INDICATORS**
- Statistical integration
- Genotype level (mostly)
- New computation for each scientific question
- One raw dataset → many computed datasets

**Data publication**

- One Data Publication by datasets.
- **Platform IS**
  - Phenomic, plant level
- **FAIR Data Repositories**
  - Reduced



Data

| Genotype | traitement | Fusariose |
|----------|------------|-----------|
| oisson   | low input  | 5         |
| oisson   |            |           |
| Charger  | low input  | 1         |
| Charger  | high N     | 2         |

Knowledge

Variety charger

intensiv cultural practice

**INRAe**

Plant science data management and integration: the heterogene

11 March 2022 / ABLS4 / Cyril Pommier

# > Plant Genetic variation

- Variability of the genotypes (AKA varieties, accessions, germplasm)

- Sequencing (GBS), Chips, …

- Raw data : reads

- Aligned data : VCF

- Paradigm: Raw data is too big, easy to generate ➔ keep only Variation

- But: realign to a new genome version, or to another reference variety ?

- ➔ Raw data can be interesting to keep too

**INRAe**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 12

# ❯ FAIR For plant science

- Phenotyping: Raw and derived data

- Genotyping: Computed data (plus option for Raw data).
  - Applies to other OMICS

- Solutions on the data lifecycle
  - Data standardisation
  - Data repositories for publication
  - Data findability / discovery

# PLANT DATA STANDARDS : WHY

**INRAe**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 14

# Why should we standardize data?

- Allow anyone (including yourself) to reuse it: metadata about the experiment (who did it, for what purpose, where and how)

- Enable data integration with other types of data: Linked data between datasets using identification of pivot objects

Phenotype 1 = measurement on a genotype in an environment-GPS1-time1

Phenotype 2 = measurement on a genotype in an environment-GPS2-time2

Genotype = observed marker's alleles on a genotype

Climate 1 = climatic data at GPS1-time1

- To enable knowledge discovery: metadata about the experiment, controlled vocabularies, ontologies

**INRAE**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 15

# PLANT DATA STANDARDS : WHO

**INRA℮**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 16

**National Networks**

**European Networks**

**Global Networks**

**International data standards**

BrAPI
Web services

RDA
RESEARCH DATA ALLIANCE

miappe
Minimal information

Crop Ontology
for agricultural data
Controlled vocabularies
Trait dictionaries

MCPD

INRAe

# Sharing standards: standards registries

**INRA℮**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

18

p. 18

# Community driven recomendations and registries

- WheatIS: http://wheatis.org/DataStandards.php

- RDM Toolkit

  ◆ **https://rdmkit.elixir-europe.org/**

- Community story



**F1000Research**
Open for Science

BROWSE    GATEWAYS & COLLECTIONS    HOW TO PUBLISH    ABOUT    Search

**RDMkit**    Data management    About    Contribute    ⊙ GitHub    🔍 Search    Check for updates

**Data management**

Data life cycle ⌄
Your role ⌄
Your domain ⌄
Your tasks ⌄
Tool assembly
National resources
All tools and resources
All training resources

**Are you working with data in the Life Sciences? Do you fee overwhelmed when you think about Research Data Manaved**

The ELIXIR Research Data Management Kit (RDMkit) is an online guide cont management practices applicable to research projects from the beginning Developed and managed by people who work every day with life science d has guidelines, information, and pointers to help you with problems throu life cycle. RDMkit supports FAIR data — Findable, Accessible, Interoperable by-design, from the first steps of data management planning to the final st data in public archives.

The RDMkit organises information into the six sections displayed below, wh interconnected but can be browsed independently.

## Data life cycle

ARTICLE

Developing data interoperability using standards: at community use case [version 2; referees: 2 ved]

Yeumo[1], Michael Alaux ⓘ[2], Elizabeth Arnaud[3], Sophie Aubin[1], Ute Baumann[4], che[5], Laurel Cooper ⓘ[6], Hanna Ćwiek-Kupczyńska[7], Robert P. Davey ⓘ[8], an Fulss[9], Clement Jonquet ⓘ[10,11], Marie-Angélique Laporte[3], Pierre Larmande ⓘ[12,13], nier ⓘ[2], Vassilis Protonotarios ⓘ[14], Carmen Reverte ⓘ[15], Rosemary Shrestha[9], rats[16], Aravind Venkatesan ⓘ[12], Alex Whan[17], ✉ Hadi Quesneville ⓘ[2]

etails

This article is included in the Global Open Data for Agriculture and Nutrition gateway.

# PLANT DATA STANDARDS : WHAT

# Data standards for FAIR

## Semantic

- Description of the data
- Controlled vocabularies: term name and definitions
- Ontologies: semantic links between terms
- *Biologist* driven

## Structure

- Formatting and Organizing the data
- Data Models
- Standards : CSV, VCF, GFF, MIAPPE (www.miappe.org) , etc...
- *Biologist & Computer scientist* driven

**Persistent Unique Identifiers**
URI, gene ID, accessions ID, Trait ID, DOI,…

## Technical

- Data integration and sharing
- Interoperability : tools and systems
  - GA4GH
  - Breeding API www.brapi.org
- *Computer scientist* driven

INRA℮

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 21

# Semantic Standard: Ontologies

- Annotating one object
  - Protein, gene
  - Plant, plant anatomy, …

- Atomic concept
  - protein function
  - cellular localization
  - …

# Semantic Standard: Ontologies for Phenotype

- Describing traits/features in specific plant species
- Crop Ontology Trait + Method + Scale Semantic model

AgroPortal LIRMM

Ontology Lookup Service
Home  Ontologies  Documentation  About

Crop Ontology
for agricultural data

Variable identification: Plant height example

**Trait**  +  **Method**  +  **Unit**

M1: Total height

M2: First tassel branch

M3: Last expanded leaf

M4: Youngest growing leaf

…There is an uncountable number of combinations…
Each trait, method and unit has to be identified if we want to share and reuse data

T1: Plant Height

M5: Highest pixel
corresponding to plant

U3: pixel

Slide from L. Cabrera-Bosquet

# Phenotype Structure Standard



## Minimal Information About Plant Phenotyping Experiment : version 1.1 (Jan 2019)

## www.miappe.org

- Many stakeholders
  - **Elixir, Emphasis, Bioversity, North American PPN**

- Open Community:
  - **Request for comments**
  - **Github Feature requests**
  - **Mailing lists**
  - **Meetings & Workgroups**

- Crops and woody plants

Papoutsoglou *et al*. (2020) Enabling reusability and interoperability of plant phenomic datasets with MIAPPE 1.1. New Phytol, 227:260-273; https://doi.org/10.1111/nph.16544

# Phenotype <u>Structure</u> Standard



Minimum Information for Biological and Biomedical Investigations

A collection of the historical MIBBI foundry reporting guidelines. The minimum information standard is a set of guidelines for reporting data derived by relevant methods in biosciences. If followed, it ensures that the data can be easily verified, analysed and clearly

- Biologist Friendly
  - Clear definitions and examples
  - Excel templates
  - Trainings

- Minimal and sufficient list of metadata:
  - The objective of the experiment
  - Who contributed to the experiment
  - What were the experimental procedures
  - What was the biological material experimented
  - …

**INRAe**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 25

# Phenotype Technical Standard, MIAPPE Implementations

- Ontology, OWL Implementation
  - https://github.com/MIAPPE/MIAPPE-ontology
  - http://agroportal.lirmm.fr/ontologies/PPEO
  - Data model representation
  - Formal concepts and constraints

- File Archive
  - ISA Tab: data + metadata
  - RO Crate studies

- Web Services
  - Breeding API
  - International collaboration
  - Standard Open Web Service API
  - Information Exchange, Main target: Breeding
  - Excellence in Breeding platform (CGIAR, Peter Selby)



**BrAPI**

# Data Integration between silos, From Phenotyping to Genotyping

Identifying key resources/pivot objects

# Data Integration between silos, From Phenotyping to Genotyping



Community data discovery portals

INRAⒺ

Plant science data management an...

11 March 2022 / ABLS4 / Cyril Pom...

# Global Data discovery portal

Dispersed data          Heterogenous data          Dedicated repositories & Archives

# FAIDARE: Global Data discovery portal

# ❯ Take Home Message

- Data integration relies on a complex lifecycle

- Both:
  - Deriving and reducing data
  - Linking different datasets

- All steps must be defined
  - → data management plan

- Not all step of data must be shared

- Raw and final data should be shared

- With sufficient provenance

- Open science (policy): Publication and Findability are keys (Data tombs effect)

**INRAƐ**

Plant science data management and integration: the heterogeneity and dispersion challenge
11 March 2022 / ABLS4 / Cyril Pommier

p. 32

# Aknowledgments